# SIFTER: Search Services for Digital Libraries

# CLIOH: Cultural digital Library Indexing Our Heritage

Mathew J. Palakal

**SIFTER Research Team**

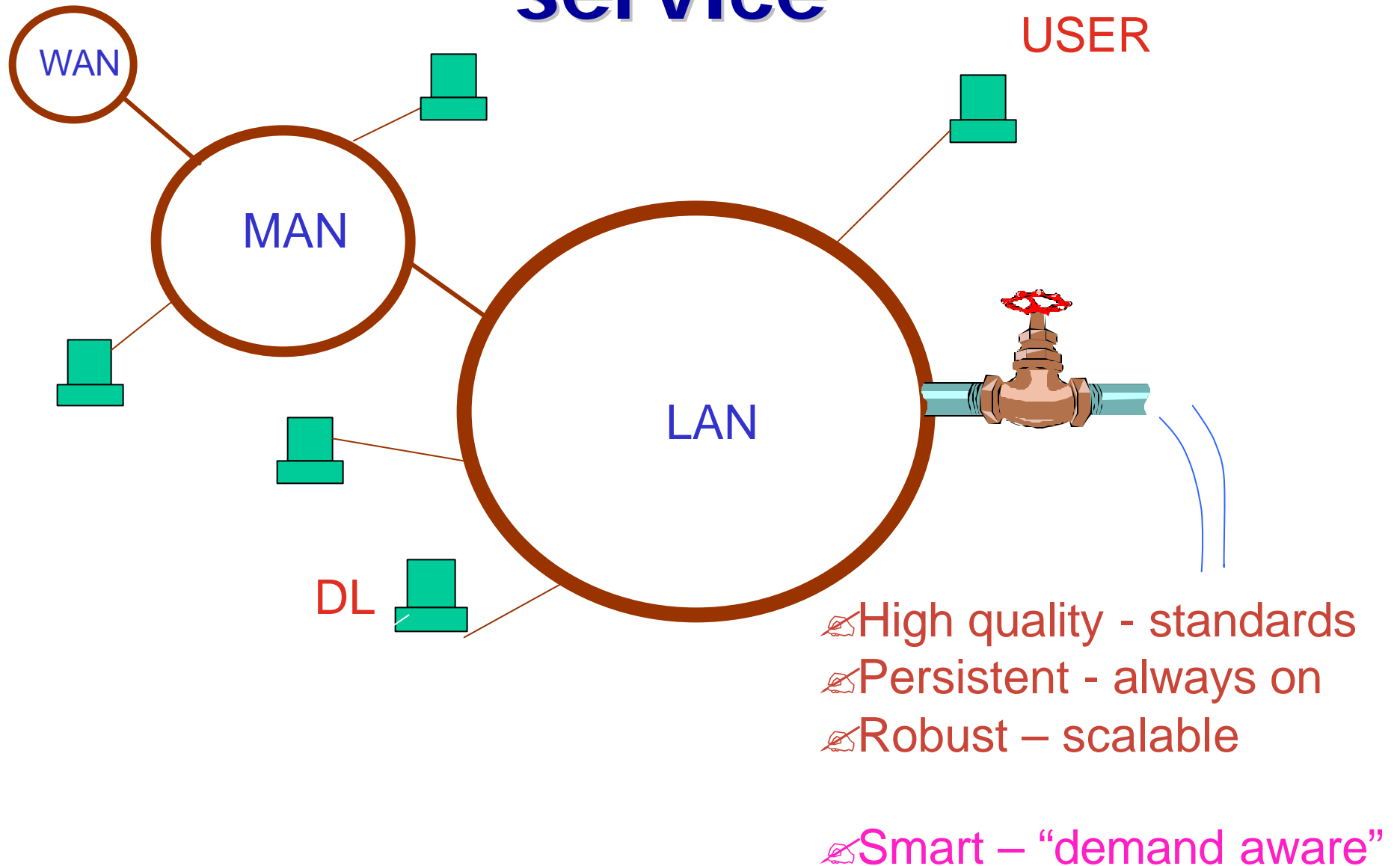**CLIOH Research Team**

Indiana University Purdue University Indianapolis
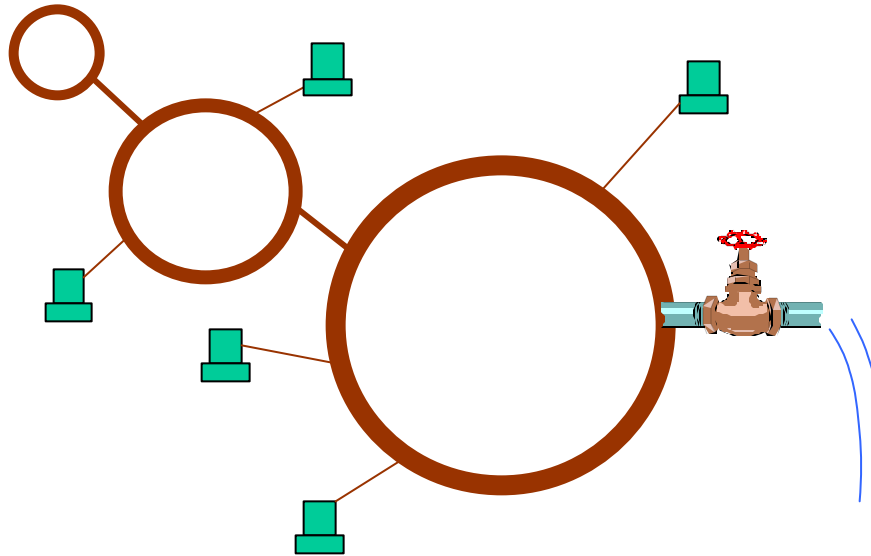Indiana University Bloomington

# Search

- ✍ **Information searching is a basic necessity …**
  - ✍ **Critical to the usefulness of a digital library**

- ✍ **Information available through a digital library may actually come from many different sources (both historical and recent)**

- ✍ **Users may need access to multiple digital libraries – distributed across the globe**

# Search as a "utility" service



WAN

MAN

LAN

USER

DL

- High quality - standards
- Persistent - always on
- Robust – scalable

- Smart – "demand aware"

# Effective Search Service



- ?Organization
  - ?Aggregation
  - ?Representation
  - ?Classification

- ?Matching
  - ?Delivery media & devices: (customization)
  - ?Users' interests: (query & profile personalization)

- ?Presentation & interaction
  - ?Prune, cluster, rank, format, visualize

# Key Challenges

1. Data Diversity
   - Diverse sources
   - Numerous formats
   - Heterogeneous content
2. Dynamic Environment
   - Content drift
   - Quality change
3. User needs
   - User's demands are context-sensitive
   - User's interest vary and may change over time

# Rising to the Challenge

**?**

SIFTER

Smart Information Filtering Technologies for Electronic Resources

Developing algorithms and systems that utilize both IR and AI approaches

# Problem 1: Data Diversity

- Need to identify document semantics: labels, terms and concepts

- Need to identify associations among concepts, terms or labels

# Data Diversity: SIFTER Solutions

- **Representation**:
  - Use of thesauri
  - Algorithms to convert data elements to efficiently computable structures

- **Classification**:
  - Use of comprehensive classification schemes
  - Algorithms to cluster or classify to higher level representations

# Problem 2: Dynamic Sources

- Local users -> local vocabularies and functions

- New "vocabularies" are introduced and they need to be discovered

# Dynamic Sources: SIFTER Solutions

- Distributed knowledge and functions using multi-agent architecture

- Vocabulary discovery based on discriminatory power

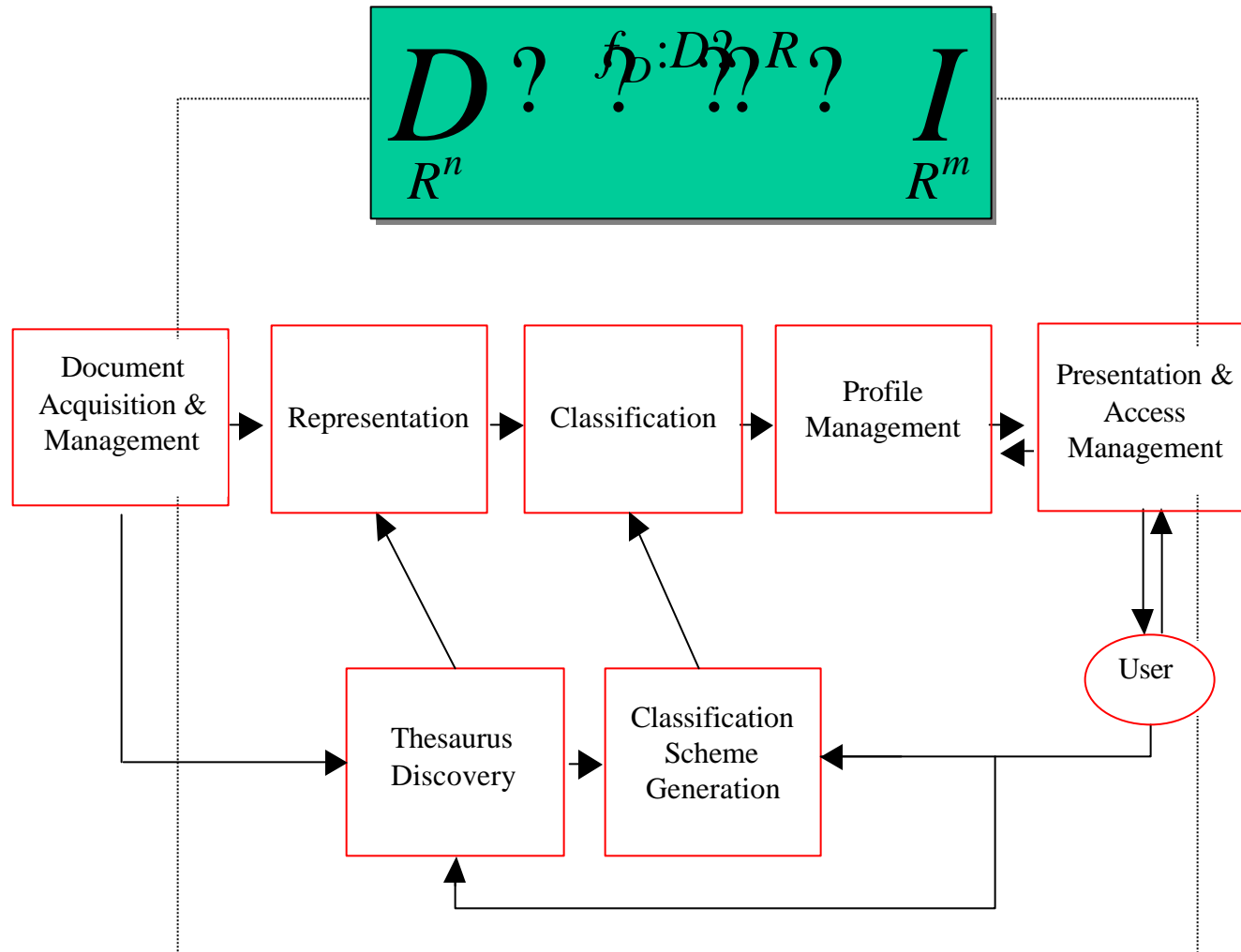- Classification scheme generation and ongoing replenishment

# Problem 3: User Need

- Interest information usually covers a small subset of the universe of topics (a profile), therefore service has to be "personalized"

- Identification and capturing of profile cannot rely directly and exclusively on the user

- Interest may not be constant over topics and can change gradually or rapidly

# User Need: SIFTER Solution

- User profile modeling: capture interest in a representation that supports topic "exploration" and reduces user involvement

- Promote control and convenient modification

- Detect "interest shifts"

- Support model and domain visualizations

# Modeling the Information Filtering Process

$$D \ ? \quad f_p{:}D? \ ?R \ ? \quad I$$

$$R^n \qquad\qquad\qquad R^m$$

| Document Acquisition & Management | Representation | Classification | Profile Management | Presentation & Access Management |
|---|---|---|---|---|

| Thesaurus Discovery | Classification Scheme Generation | User |
|---|---|---|

# Capturing User's Interest

- Explicit (topics), rating content, and user behavior

Lam, W. & Mostafa, J. "Modeling User Interest Shift Using a Bayesian Approach".
*Journal of the American Society for Information Science & Technology*, 52(5), 416-429, 2001

# Representation & Classification

- Representation:
    - Use of thesauri
    - Algorithms to convert data stream to efficiently computable structures
- Classification:
    - Algorithms to cluster or classify to higher level representations

M. Palakal, Y. Mulong, S. Mukhopadhyay, The Effects of Dynamic Clustering in Information Filtering, *Information Processing and Management*, 2002 (in preparation).

# Automated Approaches

- Learned from existing classification results – used PUBMED for training
- Developed algorithms for vocabulary and association discovery

| MeSH Classes |
| --- |
| |
| Cell Adhesion |
| Cell Communication |
| Cell Death |
| Cell Movement |
| Cell Survival |
| Endocytosis |
| Antibody Formation |
| Autoimmunity |
| Immunocompromised Host |
| Cytotoxicity Immunologic |
| Immune Tolerance |
| Immunity Cellular |
| Regeneration |
| Evolulution |
| Complement Activation |

| Automatically Produced Classes |
| --- |
| |
| Cell, Binding |
| Cell, Adhesion, Growth, Antigen |
| Communication, Death |
| Apoptois |
| Migration |
| Production, Motility |
| Tolerance |
| Virus |
| Endocytosis, Receptor |
| Antibody, Serum |
| Autoimmune |
| Tumor |
| Immunocompromised, Infected |
| Cytotoxic |
| Immune, Cell, Response, Gene, Class |
| Regeneration |
| Evolution, DNA |
| Complement, Activation, Plasma, Membrane |
| Transplant |
| Muscle |
| Expression |

Mostafa, J., & Lam, W. "Automatic Classification Using Supervised Learning in a Medical Document Filtering Application." *Information Processing & Management*, 36(3), 415-444, 2000

# Interactive Term and Cluster Discovery

# Diverse Sources: Distributed Services

**D-SIFTER**



AGENT 1    AGENT 2

SERVER

AGENT 3    AGENT 4

Distributed knowledge & Local functionality

**SIFTER-II**



users

sifter server

user agents

administrator agent

domain agents

classification agents

wrapper agents

centroid generator service

database

database

Distributed knowledge & Distributed functionality

**Raje, R., Qiao, M., Mukhopadhyay, S., Palakal, M., & J. Mostafa, J. "Homogeneous Agent-based Distributed Information Filtering",** *Cluster Computing,* **2002 (in press)**

**Raje, R., Qiao, M., Mukhopadhyay, S., Palakal, M., & J. Mostafa, SIFTER-II: A Heterogeneous Agent Society for Information Filtering,** *Proceedings of ACM Symposium on Applied Computing, SAC'01,* **pp. 121-123, Las Vegas, Nevada, 2002.**

# Evaluation of Distributed Filtering

- Progressively larger number of users -> larger agent community -> time to classify decreases
- Growing user community -> increasing number of agents -> precision suffers moderately but recall improves
- With increasing number of "user agents" classification efficiency improves



Processing Time          Filtering Performance

# SIFTER vs. the best in TREC
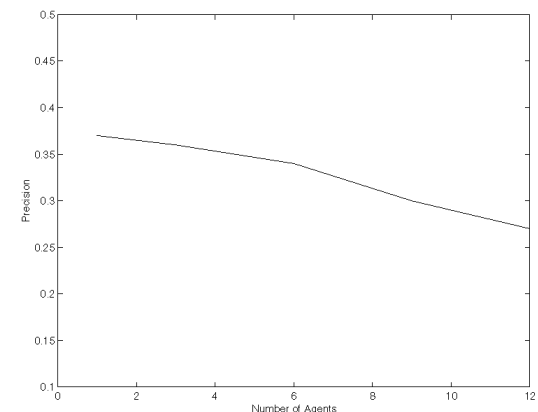
- TREC 2000 Filtering Track OHSUMED collection was used
- 293,856 documents in the test set
- 4967 topics (include OHSU and MeSH topics)
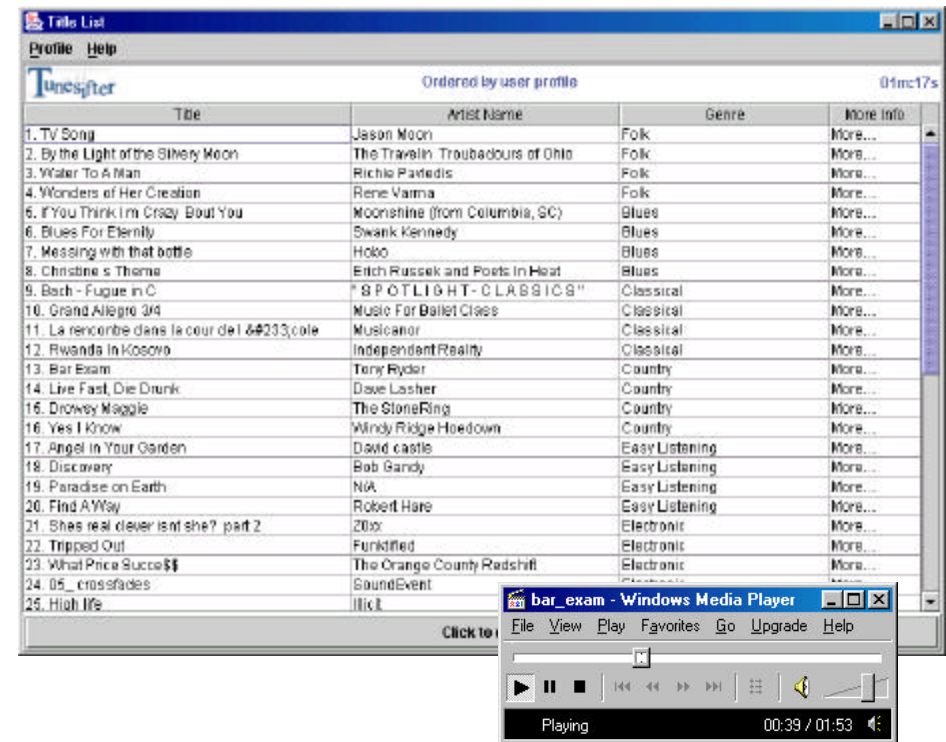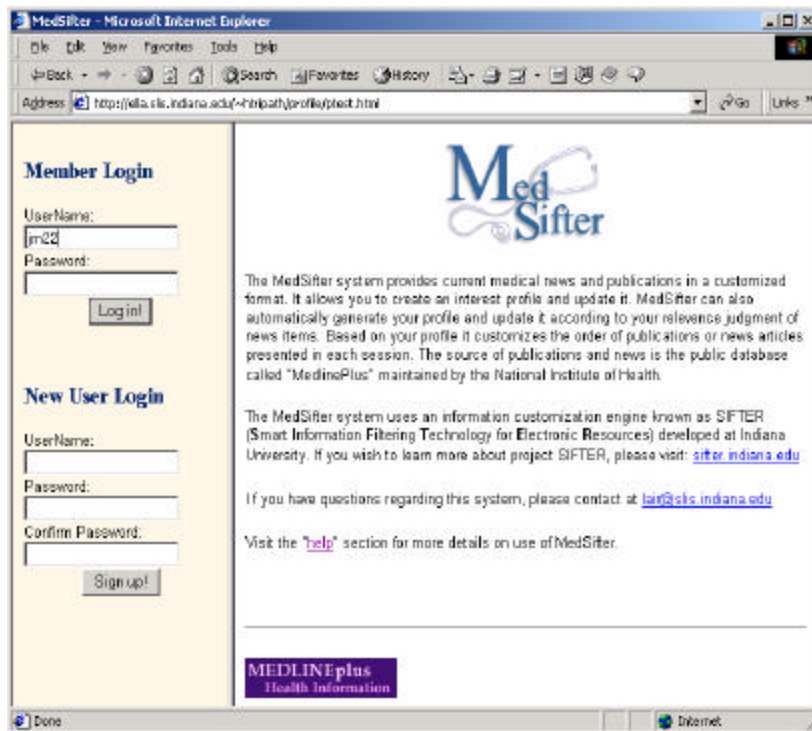- Evaluated BOTH effectiveness and efficiency

| SYSTEMS | MnT9P | MnT9U | Proc. Time/doc (in mSec) |
|---|---|---|---|
| Fudan | 31.7 | -1.1 | NA |
| Microsoft | 30.5 | -5.3 | NA |
| CMU-Y | 26.1 | -26.9 | NA |
| KAIST | 20 | 12.2 | NA |
| SIFTER (Theta 0.6) | 30.6 | -6.5 | 6165.8 |
| D-SIFTER (3 agents) | 29.9 | -8.5 | 1500.7 |
| D-SIFTER (6 agents) | 28.8 | -11.5 | 374.1 |
| D-SIFTER (9 agents) | 25.5 | -23 | 162.7 |
| D-SIFTER (12 agents) | 22.9 | -35 | 97.4 |
| SIFTER-II (3 user agents) | 29.1 | -10.5 | 1602.5 |
| SIFTER-II (6 user agents) | 27.7 | -14 | 523.2 |
| SIFTER-II (9 user agents) | 24.4 | -26.5 | 329.4 |
| SIFTER-II (12 user agents) | 22 | -38.5 | 192.1 |

S. Mukhopadhyay, S. Peng, M. Qiao, R. Raje, J. Mostafa, and M. Palakal, Distributed Multi-Agent Information Filtering, *ACM Transactions on Information Systems*, 2002 (pending review).

# Diverse Formats

- Developing systems for health news (text), scholarly research publications (text), music (audio) and cultural information (all major formats)

    - MedSIFTER

    - TuneSIFTER

    - BioSIFTER

    - ViewFinder (CLIOH)
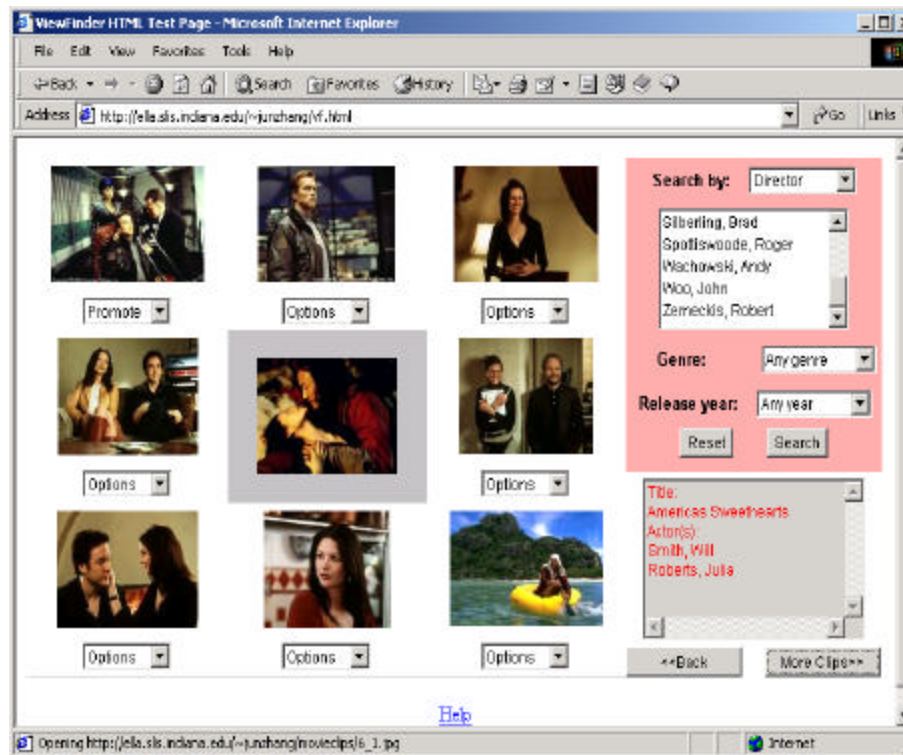
# Systems for Different Data Formats: MedSIFTER

# Systems for Different Data Formats: BioSIFTER



- Text Data
- Sequence Data
- Structural Data

M. Palakal, S. Mukhopadhyay, J. Mostafa, R. Raje, M. N'Cho, and S.K. Mishra, *An Intelligent Biological Information Management System*, *Bioinformatics*, 2002, (in press).

# Systems for Different Data Formats: CLIOH



ViewFinder

- Video Data
- Audio Data

# Beyond Current Challenges

- Cross-format, cross-language, and cross-domain information synthesis in real-time.

- Distributed DLs with both Data & Services

- Integrating Web Services with Multi-agent Searching

# Acknowledgment

This project was funded in part by:

- The National Science Foundation (DLI-II and ITR I)
- Eli Lilly & Co., Indianapolis, Indiana

## SIFTER Team:

- Mathew Palakal, Snehasis Mukhopadhyay, Javed Mostafa, Rajeev Raje
- Students: Mulong Yu, Matthew Stephens, Mingyaung Qiao, Shengquan Peng, Luz Quiroga, John Fieber, Vijay Vij